

## A Part of Speech Tagset for Sinhala

This document summarises the issues identified in the discussions held while preparing the part of speech tagset for Sinhala. The tags identified are given at end of the document.

### Issues

#### 1. Word Separation

In Sinhala, there is no generally accepted set of rules for word separation. But there seem to be groups of writers who conform to certain rules. This fact is reflected in the 10 million word Sinhala corpus too. A decision has to be taken to preserve the consistency of the part of speech tagging process.

##### 1.1. Separation of Particles and Post-positions

E.g. ගියේය / ගියේ ය, ඔහුද මමද, ඔහුගේ, ඔහු ගේ, ඒ ගැනම

##### 1.2. Separation of Compound Nouns

E.g. විශ්ව විද්‍යාලය

#### 2. Segmentation

It is the practice in Sinhala writing to combine two words belonging to certain grammatical categories to form one orthographic word.

E.g. ලස්සනයි, විවාලෙන්, විසිනි, බැවිනි, හෙයිනි, හැකිදැයි, ඇතැයි, යැයි

#### 3. Multiword

Certain word combinations/phrases can function as one grammatical category.

E.g. යටත් පිරිසෙයින්, කරණ කොට, යළි යළින්, බල බලා, නට නට

#### 4. Unknown Grammatical Categories

It is hard to determine the grammatical categories of certain words as they do not show the characteristics of any known grammatical category by themselves, particularly in the initial part of some compound verbs.

E.g. මිය ගියේය, අසු වී, පත් වූහ, අවදිවූයේ, පහළ වීමට

#### 5. Membership criteria for post-positions and particles of *Nipatha*

By definition a post-position follows a noun or a noun phrase. In traditional grammar of Sinhala called them as *nipatha*. But in some cases *nipatha* follow too. In such context it is more appropriate to label them as particles.

E.g. පොත බලන තෙක් / විට / කල්හි / අයුරු / නිසා  
මගේ පොත ගැන / තුළ / නිසා

#### 6. Pronouns

One part of speech label can be assigned to all the pronouns. A detailed tag can be assigned simple table look-up by maintaining a list of all the pronouns as they form closed class.

<b>TAG</b>	<b>Description</b>	<b>Example</b>
<b>NNR</b>	Common Noun Root	මිනිස්, පුටු, බලු
<b>NNM</b>	Common Noun Masculine	මිනිසා, බල්ලා, ශිෂ්‍යයෝ, එච්චන්
<b>NNF</b>	Common Noun Feminine	නිලියෝ , ඇතින්ත
<b>NNN</b>	Common Noun Neuter	පුටු, ගම
<b>NNPA</b>	Proper Noun Animate	රත්නායකගේ
<b>NNPI</b>	Proper Noun Inanimate	රත්නපුරේ, ලක්ස්ප්‍රේ
<b>PRPM</b>	Pronoun Masculine	ඔහු, ඒකා, කොයිකා
<b>PRPF</b>	Pronoun Feminine	ඇය, ඕකි
<b>PRPN</b>	Pronoun Neuter	එය, ඕක
<b>PRPC</b>	Pronoun Common	මම, ඔවුහු
<b>QFNUM</b>	Number Quantifier	එක, දෙවනි
<b>DET</b>	Determiner	මේ, ඒ, අර, ඔය, බොහෝ, සියලු,
<b>JJ</b>	Adjective	රළු, සුමුදු
<b>RB</b>	Adverb	වහා, ඉතින්
<b>RP</b>	Particle	ම, ලු, ය, ද, නම්, වැනි, වූකලී
<b>VFM</b>	Verb Finite Main	බලයි, බැලූහ
<b>VNF</b>	Verb Non Finite	බලා, බලමින්, බලද්දී, බලනොත්
<b>VP1</b>	Verb Participle 1	බලන,බැලූ, කළ
<b>VP2</b>	Verb Participle 2	බලනු, බැලූවා
<b>VP3</b>	Verb Participle 3	බැලිය
<b>VP4</b>	Verb Participle 4	බලන්නේ, බැලූවේ, කළේ
<b>VNN</b>	Verbal Non Finite Noun	බැලීමි, බැලීලි, බැලුම්
<b>POST</b>	Postpositions	ගැන, ලෙස, සඳහා
<b>CC</b>	Conjunctions	සහ , ද, සමඟ, හෝ
<b>NVB</b>	Noun in Kriya Mula	පාඩම් කරනවා, භාජනය කරනවා
<b>JVB</b>	Adjective in Kriya Mula	කිකරු වෙනවා, එකඟ වෙනවා, අඩු කරනවා
<b>UH</b>	Interjection	අහෝ , චී, ඕක්, ආූ
<b>FRW</b>	Foreign Word	Computer
<b>SYM</b>	Not Classified	A4